

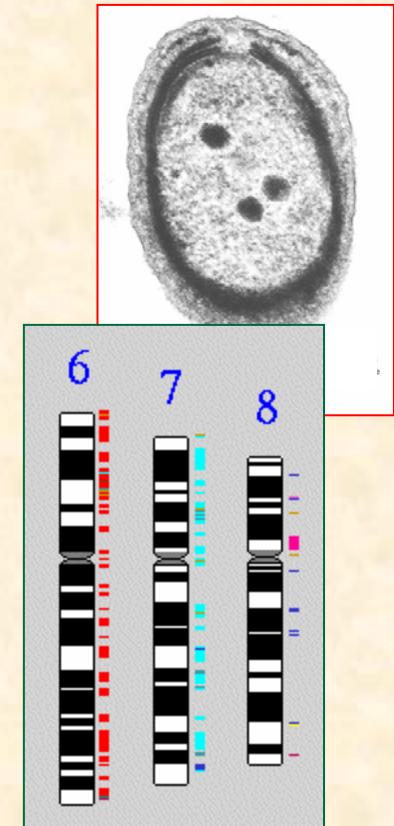
# **Genome annotation: Bioinformatics for high- throughput genomics and beyond**

**Frank Larimer**  
**Oak Ridge National Laboratory**

**9th DOE Contractor and Grantee Workshop**  
**Oakland, CA, 28 January 2002**

# Comparative Genomics Resource

- Microbial
  - JGI & public finished genomes (~80)
  - JGI & Collaborative draft genomes (~25)
- Human
  - Finished NT sequences
  - Draft Genbank entries
  - Draft Human Genome Builds
  - JGI chromosomes - 5 - 16 - 19
- Mouse and model organisms
  - Finished mouse clones / NTs
  - Mouse draft as available
  - JGI model organisms



# Community Tools and Services

- GRAIL EXP **eukaryotic genome analysis**
- Microbial Analysis **Generation, Critica, Glimmer**
- HPC Blast, Pfam, InterPro & other Tools
- Genome Analysis Pipeline **human and microbial**
- PROSPECT **protein structure prediction server**
- SRS, GDB, Swissprot, HOBACGEN **database servers**
- Genome Channel / Catalog **eukaryotic and microbial genomes / draft chromosomes**

# Web access to annotation

Computational Biology at ORNL

Home   About Us   Analysis Tools   Information Resources   Projects & Research

Channel • Generation • Grail • GrailEXP • Pipeline • Parser • PROSPECT



## Genome Channel

Search Genome Channel for:  limit search to [organism](#):

Enter keywords or phrases (examples: AC005400 or "cytokeratin 8") Tip: put an exact phrase in quotes

**Eukaryotes**  
Finished and draft clones and assemblies from GenBank, analyzed, updated, and summarized.

<b>Human</b> <input type="button" value="Pick one"/>	<b>Mouse</b> <input type="button" value="Pick one"/>	<b>Other Eukaryotes</b> <input type="button" value="Pick an Organism"/>
---	---	--

**Archaea and Bacteria**  
Finished microbial sequences from GenBank, analyzed, updated, and summarized.

<b>Finished Archaea</b> <input type="button" value="Pick an Organism"/>	<b>Finished Eubacteria</b> <input type="button" value="Pick an Organism"/>
--	---

**Other Microbials**  
Draft microbial genomes sequenced by the DOE [Joint Genome Institute](#) (JGI) and other centers, analyzed, updated and summarized. These data are a work in progress and are subject to frequent change.

**JGI Draft Microbial**

**Services**   **Data**   **Help**  
[Keyword Search](#)   [Summary](#)   [Usage](#)  
[Nucleotide Search](#)   [Download](#)   [Credits](#)  
[Protein Search](#)

---

[SiteMap](#)   [Feedback](#)   [Life Sciences Division](#)   [ORNL](#)   [Disclaimer](#)   [Webmaster](#)

OAK RIDGE NATIONAL LABORATORY  
U. S. DEPARTMENT OF ENERGY



# Web access to annotation

Computational Biology at ORNL

The screenshot shows the "Genome Channel" section of the website. At the top, there's a navigation bar with links to Home, About Us, Channel, and other sections like Analysis Tools, Information Resources, and Projects & Research. Below the navigation is a search bar labeled "Search This Organism for:" with a text input field and a "GO" button. To the right of the search bar is a "Human Overview" section with a "Links" list containing numbered links from 1 to 22, with "Unknown X Y" also listed. Below the links is a section titled "GenBank Contigs (Accession Number format NT\_xxxxxx)" with a bulleted list of details about the contig assembly. On the left side of the main content area, there's a sidebar with sections for "Search Genome Chan", "Enter keywords or phrases", "Finished and draft clones and Human", "Human Data Sets", "HTML or ASCII", "Draft microbial genomes sequ", and "JGI Draft Microbial". Each of these sections has a "Pick one" dropdown menu. At the bottom left of the sidebar is a "SiteMap" link.

# Web access to annotation

## Computational Biology at ORNL

The screenshot displays the "Genome Channel" section of the Computational Biology at ORNL website. The top navigation bar includes links for Home, About Us, Channel, Genome, and Analysis Tools, Information Resources, Projects & Research.

**Search This Organism for:**  Enter keywords or phrases (examples: AC005400 or "cytokeratin 8") Tip: put an exact phrase in quotes

**Links**

- Cl
- Ur

**GenBank** 19p13.3  
• M   
• Es   
• 29   
• As   
• 11   
• 74   
• 92

**Chromosome Overview**

GO to another organism   
Pick an Organism

**Summary**

59 Contigs (6.79 MB), 794 Clones (75.97 MB), 755 GenBank CDSs, 2284 Genscan Genes (1590 with BLAST hits), 2998 GrailEXP Genes (1711 with BLAST hits and 2376 with EST evidence)

**Links to Contigs**

Contig	Begin Position	End Position	Length	Band	Color
NT_025171.1	1	43299	43299	p13.3	
NT_025146.2	93300	216761	123462	p13.3	
NT_011287.1	266762	339739	72978	p13.3	
NT_011269.3	389740	889240	499501	p13.3	
NT_025196.2	939241	1079045	139805	p13.3	
NT_011284.2	1129046	1255132	126087	p13.3	
NT_011227.3	1305133	1532235	227103	p13.3	
NT_011286.1	1582236	1655184	72949	p13.3	
NT_011286.2	1706185	2126000	4154144	p13.3	

# Web access to annotation

Computational Biology at ORNL

Home About Us  
Channel Channel

Search Genome Chan

Enter keywords or phrases

Finished and draft clones and Human

Pick one

Finished microbial sequences

Finished Archaea

Pick an Organism

Draft microbial genomes sequ a work in progress and are sub JGI Draft Microbial

Pick an Organism

[SiteMap](#)



Computational Biology at ORNL



Search This

Organism for:

Enter keywords or phrases (examples: AC005400 or "cytokeratin 8") Tip: put an exact phrase in quotes

GO

Human Data Sets

HTML or ASCII

Many links may take a long time to finish.

Contigs [HA](#)

Clones [HA](#)

GenBank CDSs [HA](#)

GrailEXP Gene [HA](#)

Models [HA](#)

GrailEXP PolyAs [HA](#)

Genscan Gene [HA](#)

Models [HA](#)

Genscan PolyAs [HA](#)

tRNA Genes [HA](#)

STS [HA](#)

Repeats [HA](#)

CpGs [HA](#)

Links

• Cl

Ur

GenBa

• M

• Es

• 29

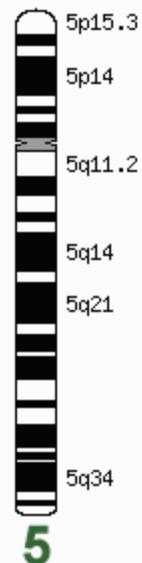
• As

• 11

• 74

• 92

more

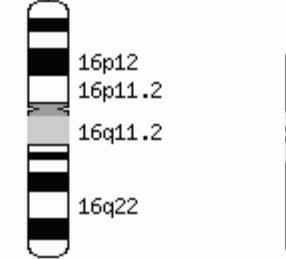


5

## Computational Annotation of Chromosomes



JOINT GENOME INSTITUTE



16



19

Computational Analysis by



# Web access to annotation

Computational Biology at ORNL

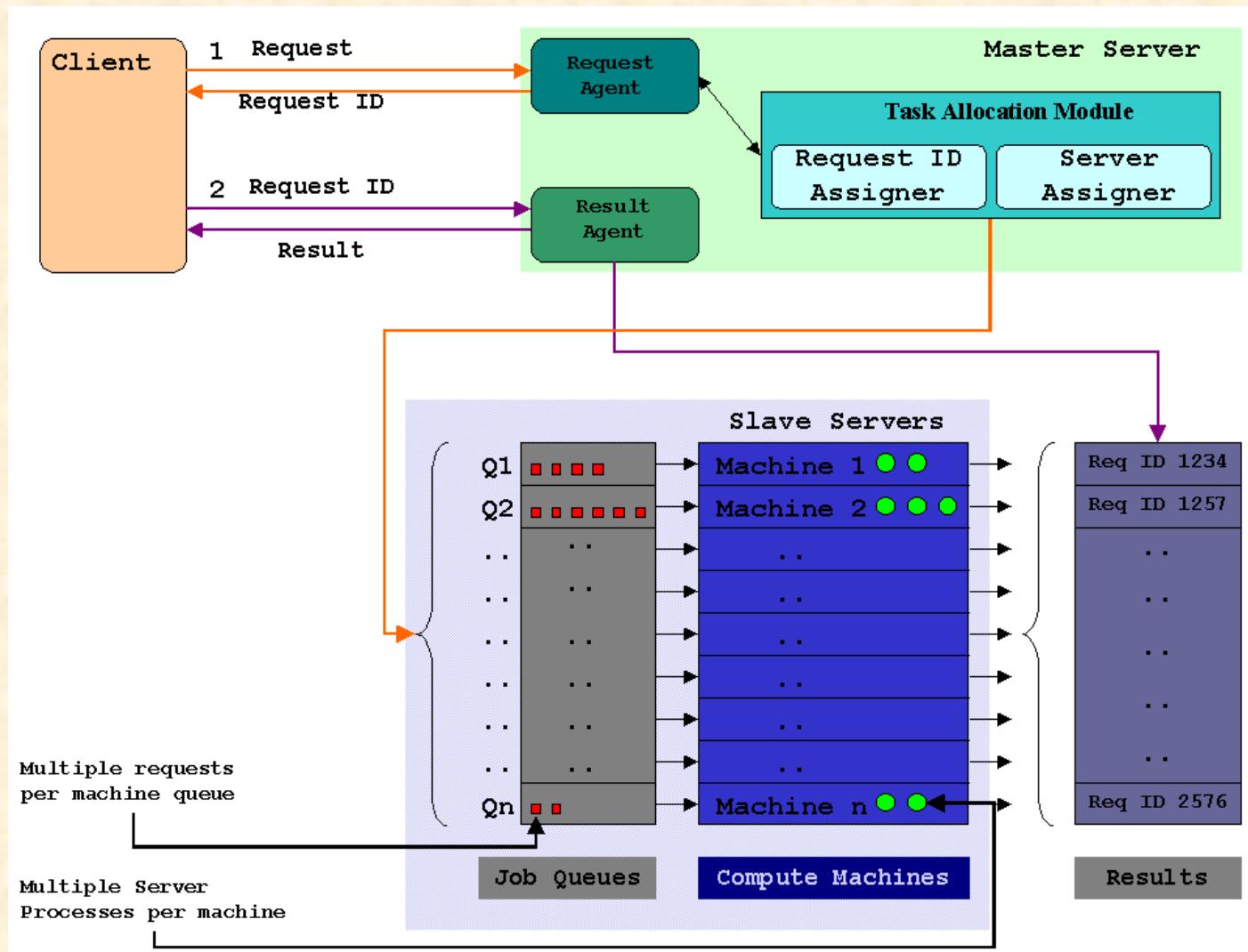
The screenshot displays a web-based genome annotation tool. At the top, there are three header sections: "Home", "About Us", "Channel", and "Genome Channel". Below these, a search bar asks "Search This Organism for:" followed by a text input field containing "Enter keywords or phrases (examples: AC005400 or "cytokeratin 8") Tip: put an exact phrase in quotes". To the right, another search bar asks "Search This Organism for:" followed by "GO to all". A "Links" section includes "Human Data Sets" with options like "HTML or ASCII", "GenBank", "Clones", "Gene", "Exon", "EstRef", "Gene", "Models", "PolyAs", "Genscan Gene", "Models", "Genscan PolyAs", "tRNA Genes", "STS", "Repeats", and "CpGs". The main content area features a "Chromosome Overview" for "Human / Human / Chromosome 19". It shows a chromosome ideogram for chromosome 19 with bands 5p15.3, 5p14, and 5q11.2. Below the ideogram is a "Sequence" track labeled "hv". A large green box highlights the "Computational Annotation of Chromosomes" section, which includes a "GenBa" interface with a "File" menu (File, View, Features, Windows) and a zoom scale (2.0, 7.0, 6028 base). The sequence track shows various genomic features and annotations, including a red box highlighting "EstRef 1 (3248..6052), 3 est hits 5 est AA292232E (6009..6052)". A detailed table of repeats is shown in the bottom left, and a sequence viewer is in the bottom right.

Name	Strand	Start	Stop	Length
1	f	3236	6085	2850
2	f	14771	15197	427
3	f	16746	17235	490
4	f	18704	19806	1103
5	r	450	6751	6302
6	r	7441	7602	162
84	r	11203	11456	254

OAK RIDGE NATIONAL LABORATORY  
U. S. DEPARTMENT OF ENERGY

UT-BATTELLE

# Web-accessible Analysis Pipeline



# Infrastructure Detail

## *The High-performance Computing Environment*

- Hardware:
  - Eagle:IBM SP3 184 nodes (736 Power 3+ cpus), ~1 Tflops
  - Falcon:Compaq 64 nodes (256 Alpha cpus), ~300 Gflops
  - Colt :Compaq 16 nodes (64 Alpha cpus), ~100 Gflops
  - Htorc : Intel/Linux 64 nodes (128 Pentium III cpus)
- Storage:
  - HPSS (High Performance System),
  - up to 100TBytes.
- Applications:
  - GIST (Genome Integrated Supercomputing Toolkit) , includes MPP NCBI BLAST, MPP HmmerPfam, GrailEXP, PROSPECT, InterPro, genome assembly



# JGI Microbial Sequencing

- Final annotation:
  - *Prochlorococcus marinus* MED4
  - *Rhodopseudomonas palustris*
  - *Synechococcus* WH8102
  - *Prochlorococcus marinus* MIT9313
  - *Nitrosomonas europaea*
- Drafted, finishing:
  - *Nostoc punctiforme* - 198 contigs
  - *Cytophaga hutchinsonii* - 98 contigs
- Microbe Month: 15 species sequenced to deep draft (6-8x coverage)
- Draft output
  - ~93 million bp analyzed
  - ~85,000 gene models
- Microbe drafting ongoing in 2002...

# DOE JOINT GENOME INSTITUTE

Operated by The University of California for The US Department of Energy

## News and Events

### Who We Are

- Mission Statement
- JGI Members/Partners
- Production Genomics Facility
- Organizational Chart
- Contact Us / Employment
- FAQs

### JGI Programs

- Human Genome Project
- Microbial Genomics
- HSA19/mouse Comparative Sequencing
- Fugu Genome Project
- Ciona Genome Project
- White Rot Genome Project
- Comparative Genomics
- Instrumentation
- Computational Genomics
- Research and Development

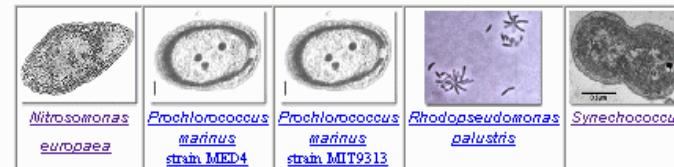
### JGI Home

- Production Statistics
- Production Protocols
- JGI Internal Site

### Download JGI Sequences

### Genome Links

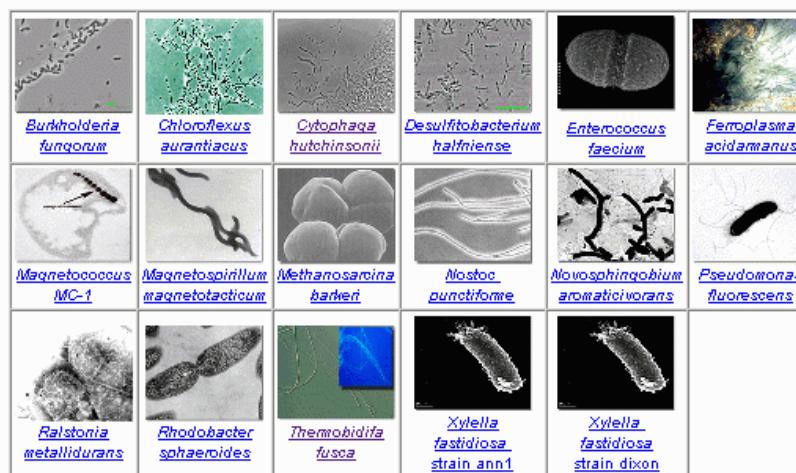
## Completed Microbial Genomes



## 2002 Microbial Genomes



## 2001 Microbial Genomes





Produced  
for the [Joint](#)  
[Genome](#)  
[Institute](#)  
Microbial  
Sequencing  
program

DOE  
[Microbial](#)  
[Genomes](#)

# DRAFT \* DRAFT \* DRAFT

## *Prochlorococcus marinus sp. MED4* analysis files

Version 01sep00

1,657,995 bp in 1 contigs of 20 reads or greater, 31 % GC  
1694 candidate protein-encoding gene models

<b>Modeled Genes Organized by:</b>	<b>FASTA Files</b> <ul style="list-style-type: none"><li>• <a href="#">Sequence Data</a> (used for gene modeling)</li><li>• <a href="#">CDS models</a></li><li>• <a href="#">CDS translations</a></li></ul>
<b>Tab-Delimited Files</b> <ul style="list-style-type: none"><li>• <a href="#">Gene model summary detail tab-delimited</a></li><li>• <a href="#">NR blast top hit tab-delimited</a></li><li>• <a href="#">COGs tab-delimited</a></li><li>• <a href="#">Pfam tab-delimited</a></li><li>• <a href="#">EC top hit tab-delimited</a></li><li>• <a href="#">GC content, GC and AT skew, per gene</a></li><li>• <a href="#">KEGG Genes database top hit</a></li></ul>	<b>RNAs</b> <ul style="list-style-type: none"><li>• <a href="#">tRNA content</a></li><li>• <a href="#">16S rRNA</a></li><li>• <a href="#">23S rRNA</a></li></ul>
<b>Search the Sequence</b> <ul style="list-style-type: none"><li>• Blast <a href="#">search engine</a></li></ul>	<b>Repeats</b> <ul style="list-style-type: none"><li>• <a href="#">Summary of Repeats</a></li></ul>



Produced  
for the [Joint](#)  
[Genome](#)  
[Institute](#)  
Microbial  
Sequencing  
program

## DRAFT \* DRAFT \* DRAFT

### *Prochlorococcus marinus sp. MED4*

### files

Version 01sep00

1,657,995 bp in 1 contigs of 20 reads or greater, 31 % GC  
1694 candidate protein-encoding gene models

Modeled Genes Organized by:	FASTA Files
<ul style="list-style-type: none"> <li><a href="#">Metabolic Pathway</a></li> <li><a href="#">Functional Category</a></li> <li><a href="#">Contig</a></li> <li><a href="#">COGs Functional Group</a></li> </ul>	<ul style="list-style-type: none"> <li><a href="#">Sequence Data</a> (used for gene model)</li> <li><a href="#">CDS models</a></li> <li><a href="#">CDS translations</a></li> </ul>
Tab-Delimited Files	RNAs
<ul style="list-style-type: none"> <li><a href="#">Gene model summary detail tab-delimited</a></li> <li><a href="#">NR blast top hit tab-delimited</a></li> <li><a href="#">COGs tab-delimited</a></li> <li><a href="#">Pfam tab-delimited</a></li> <li><a href="#">EC top hit tab-delimited</a></li> <li><a href="#">GC content, GC and AT skew, per gene</a></li> <li><a href="#">KEGG Genes database top hit</a></li> </ul>	<ul style="list-style-type: none"> <li><a href="#">tRNA content</a></li> <li><a href="#">16S rRNA</a></li> <li><a href="#">23S rRNA</a></li> </ul>
Search the Sequence	Repeats
<ul style="list-style-type: none"> <li>Blast <a href="#">search engine</a></li> </ul>	<ul style="list-style-type: none"> <li><a href="#">Summary of Repeats</a></li> </ul>

DOE

Netscape

File Edit View Go Communicator Help

***Prochlorococcus marinus* MED4 genes g**

**Photosystem I**

Contig26 [126](#)- slr1834 - d1018170 - P700 apoprotein subunit V

Contig26 [127](#)- slr1835 - d1018171 - P700 apoprotein subunit V

Contig26 [61](#)- ss10563 - d1010687 - photosystem I subunit V

Contig26 [82](#)- slr0737 - d1017421 - photosystem I subunit I

Contig26 [539](#)- ssr2831 - d1019116 - photosystem I subunit I

Contig26 [434](#)- s110819 - d1018841 - photosystem I subunit I

Contig26 [435](#)- sml0008 - d1018840 - photosystem I subunit I

Contig24 [160](#)- ssr0390 - d1017452 - photosystem I subunit I

Contig26 [129](#)- slr1655 - d1019506 - photosystem I subunit I

**Photosystem II**

Contig26 [1081](#)- s110247 - d1018681 - chlorophyll a/b binding protein

Contig26 [833](#)- s111867 - d1018963 - photosystem II D1 protein

Contig26 [551](#)- slr0906 - d1011109 - photosystem II CP47 protein

Contig25 [192](#)- s110851 - d1018532 - photosystem II CP43 protein

Contig25 [191](#)- s110849 - d1018533 - photosystem II D2 protein

Contig26 [870](#)- ssr3451 - d1017825 - cytochrome b559 a subunit

Contig26 [871](#)- smr0006 - d1017826 - cytochrome b559 b subunit

Contig26 [567](#)- slr1280 - d1019017 - NADH dehydrogenase subunit

Contig26 [590](#)- ss12598 - d1018362 - photosystem II PsbH protein

Contig26 [847](#)- sml0001 - d1010736 - photosystem II PsbI protein

Contig26 [873](#)- smr0008 - d1017828 - photosystem II PsbJ protein

Contig26 [858](#)- sml0005 - d1017808 - photosystem II PsbK protein

Contig26 [872](#)- smr0007 - d1017827 - photosystem II PsbL protein

Contig26 [878](#)- sml0003 - d1017441 - photosystem II PsbM protein

Contig26 [846](#)- smr0009 - d1018363 - photosystem II PsbN protein

Contig26 [604](#)- s110427 - d1019207 - photosystem II manganese protein

Contig26 [552](#)- smr0001 - d1010784 - photosystem II PsbT protein

Contig24 [175](#)- s111398 - d1017351 - photosystem II 13 kD protein

Contig26 [413](#)- slr1645 - d1019195 - photosystem II 11 kD protein

**Phycobilisome**

Contig26 [875](#)- s111577 - d1017965 - class III phycoerythrin

**Soluble electron carriers**

Contig23 [110](#)- s111245 - d1018905 - cytochrome CytM (cytM)

Contig25 [101](#)- s110248 - d1018680 - flavodoxin (isiB)

Contig26 [378](#)- s110199 - d1010878 - plastocyanin (petE)

Contig26 [624](#)- s110100 - d1010879 - plastocyanin (petE)

Document Done



Produced  
for the [Joint](#)  
[Genome](#)  
[Institute](#)  
Microbial  
Sequencing  
program

# DRAFT \* DRAFT \* DRAFT

## *Prochlorococcus marinus sp. MED4*

### files

Version 01sep00

1,657,995 bp in 1 contigs of 20 reads or greater, 31 % GC  
1694 candidate protein-encoding gene models

Modeled Genes Organized by:	FASTA Files
<ul style="list-style-type: none"> <li><a href="#">Metabolic Pathway</a></li> <li><a href="#">Functional Category</a></li> <li><a href="#">Cluster</a></li> <li><a href="#">Draft pmar Data - Netscape</a></li> <li><a href="#">File Edit View Go Communicator Help</a></li> </ul>	<ul style="list-style-type: none"> <li><a href="#">Sequence Data</a> (used for gene model)</li> <li><a href="#">CDS models</a></li> <li><a href="#">FASTA files</a></li> </ul>

DOE

***Prochlorococcus marinus MED4 genes***

**Photosystem I**

Contig26	<a href="#">126</a>	-	slr1834	-	d1018170	-	P700 apoprotein subunit
Contig26	<a href="#">127</a>	-	slr1835	-	d1018171	-	P700 apoprotein subunit
Contig26	<a href="#">61</a>	-	ss10563	-	d1010687	-	photosystem I subunit V
Contig26	<a href="#">82</a>	-	slr0737	-	d1017421	-	photosystem I subunit I
Contig26	<a href="#">539</a>	-	ssr2831	-	d1019116	-	photosystem I subunit
Contig26	<a href="#">434</a>	-	s110819	-	d1018841	-	photosystem I subunit
Contig26	<a href="#">435</a>	-	sm10008	-	d1018840	-	photosystem I subunit
Contig24	<a href="#">160</a>	-	ssr0390	-	d1017452	-	photosystem I subunit
Contig26	<a href="#">129</a>	-	slr1655	-	d1019506	-	photosystem I subunit

**Photosystem II**

Contig26	<a href="#">1081</a>	-	s110247	-	d1018681	-	chlorophyll a/b binding protein
Contig26	<a href="#">922</a>	-	s111967	-	d1018062	-	photosystem II D1 protein

Tab	Gene	Strand	Generation	Glimmer	Critica	Stop	Best Hit	Enzyme Maps
	<a href="#">193</a>	r	202304(ATG)..200775	202304(ATG)..200775	202304(ATG)..200775	200775 (TGA)	>gi 2499981 sp Q55860 CBP1 SYNTHASE >gi 1001778 dbj BAA10617  (D64004) cobyric acid synthase [Synechocystis sp.] e-118	
	<a href="#">109</a>	f	202371(GTG)..202997	202371(GTG)..202997	202371(GTG)..202997	202997 (TAA)	>gi 1653512 dbj BAA18425  (D90914) MAF [Synechocystis sp.] 2e-31	
	<a href="#">192</a>	r	204422(GTG)..203040	204449(TTG)..203040	204422(GTG)..203040	203040 (TAA)	>gi 227610 prt 1707315B photosystem II CP43 protein [Synechococcus sp.] 0.0	
	<a href="#">191</a>	r	205482(ATG)..204406	205482(ATG)..204406	205482(ATG)..204406	204406 (TAA)	>gi 131299 sp P11005 PSBD SYNTHETIC PHOTOSYSTEM II D2 PROTEIN (PHOTOSYSTEM Q(A) PROTEIN) >gi 79676 prt JU0321 photosystem II protein D2 - Synechococcus sp. (PCC 7942) >gi 154586 (M20815) D2 thylakoid protein [Synechococcus sp.] >gi 154588 (M20814) psbD1 thylakoid protein [Synechococcus sp.] 0.0	
	<a href="#">110</a>	f	205667(ATG)..206224	205667(ATG)..206224	205667(ATG)..206224	206224 (TAA)	>gi 1723362 sp P51220 YCF4 PORPHYRINOPROTEIN HYPOTHETICAL 21.2 KD PROTEIN YCF4 (ORF186) >gi 2147565 prt S73141 hypothetical protein 4 - Porphyra purpurea chloroplast >gi 1276686 (U38804) hypothetical chloroplast ORF 4. [Porphyra purpurea] 2e-34	
	<a href="#">190</a>	r	207401(ATG)..206877	207401(ATG)..206877	207401(ATG)..206877	206877 (TAA)	>gi 1708470 sp P51230 ILVH PORPHYRINOPROTEIN ACETOLACTATE SYNTHASE SMALL SUBUNIT (AHAS) (ACETOHYDROXY-ACID SYNTHASE SMALL SUBUNIT) (ALS) >gi 2147069 prt S73151 acetylhydroxyacid synthase small chain - Porphyra purpurea chloroplast >gi 1276696 (U38804) acetylhydroxyacid synthase small subunit [Porphyra purpurea] 4e-56	4.1.3.18 - map00290 map00650 map00660 map00770
	<a href="#">111</a>	f	206398(GTG)..206880	206347(TTG)..206880	206398(GTG)..206880	206880 (TAA)	>gi 3184556 (AF052290) peptidyl-prolyl cis-trans isomerase B [Synechococcus PCC7002] 5e-18	



# Computational Biology at ORNL

[Home](#) [About Us](#) [Analysis Tools](#) [Information Resources](#) [Projects & Research](#)  
 Channel • Generation • Grail • GrailEXP • Pipeline • Parser • PROSPECT

Produced  
for the [Joint](#)  
[Genome](#)  
[Institute](#)

Microbial  
Sequencing  
program

*Prochlo*

## Modeled Genes Organization

- [Metabolic Pathway](#)
- [Functional Category](#)
- [Cluster](#)
- [Draft pmar Data - Netscape](#)
- [File Edit View Go Communicator Help](#)

Tab

- 
- 
- 
- 
- 
- 
- 

Search

OAK  
U. S. I

**Prochlorococcus marinus files - Netscape**  
 File Edit View Go Communicator Help

Version 10dec99 - Contig26 Gene 875

Gene Finders

Strand = r  
 Stop Location = 279721  
 Stop Codon = TAG

Gene Modeler	Start Location	Start Codon
Generation	279194	ATG
Glimmer	279194	ATG
Critica	279194	ATG

Blast against NR

Best Hit	Eval	Percent Identity
gi 1669716 emb CAA93118  (Z68890) class III phycocerythrin beta-subunit [Prochlorococcus marinus] gg 5734498 emb CAB52705.1  (AJ001230) phycocerythrin type III beta subunit [Prochlorococcus marinus]	9e-29	39

Refresh [BLAST](#) results

## PFam Model Comparison

Model	Description	Eval	Score
Phycobilisome	Phycobilisome proteins	9.4e-29	103.1

Refresh [PFAM](#) results

## Best NR Hit with an EC Number

No EC Pathway hits

## Synechocystis PCC6803 Comparison

Synecho ID	Description	Eval	Percent Identity	Category
/	phycocyanin b subunit	4e-18	31.617647	

## MRNA

```
ATGACAGTTCAAAAGAGTAATCAAATTATCAATGATAGAGATTAGAAAATAAAGTAATAAAAATTGAAGACAT
AAAAGAATTTAAATAACTGCAAACTCAAGATTAGTCAATAGATTCAATAACAAATAATTCTCACGCAATTGCCGTG
ATGCTGTGACTGCAATGATTGTGAAAATCAAGATTCAAGTTAACAAAAATATCTTAGATACCAACAAAGATGCT
GTTTGTCTAAGAGATGGAGAAATAATTAAAGGATTGTTCTACCTTTGATTTCTGATGACGAATCAGTTTATCTAA
AAACTGTTAAAGGATCTAAATACTTATGGCCCTTGGGTACCTCTGAAAAATGCTATCCGAGTTTTGAATTGA
TGAGAGATGCAACGATTCTGATTTAAAGTCACTGTAATTCTATGAAAGGAGAAAAGAATTCTTCTGATTTAATT
TCTAATACAGAGTTCAATTGAGAGAATAATTAACTTTAAGATAG
```

## PROTEIN

```
MTVSKSNQILSNRDLENISNKNIEDIKEFINTANSRLDAIDSITNNSHIAAADAVTAMICENQDSVNTKISLDTTNKMS
VCLRDGEIILRIVSYLLISDDESVLSKNCLKDLKNTYLALGVPLKNAIRVFELMRDATALSDLKSTVNNSMKGEKEFLSDL
SNTEOFERIINLLR*
```

Document Done

OE

Help

## *s marinus* MED4 genes g

834 - d1018170 - P700 apoprotein subunit  
 835 - d1018171 - P700 apoprotein subunit  
 63 - d1010687 - photosystem I subunit V  
 37 - d1017421 - photosystem I subunit I  
 831 - d1019116 - photosystem I subunit  
 819 - d1018841 - photosystem I subunit  
 008 - d1018840 - photosystem I subunit  
 390 - d1017452 - photosystem I subunit  
 655 - d1019506 - photosystem I subunit

0247 - d1018681 - chlorophyll a/b binding protein  
 967 - d1018062 - photosystem II D1 protein  
 CP47 protein  
 CP43 protein

Hit	Enzyme Maps
SYRIC ACID SYNTHASE ic acid synthase [Synechocystis sp.]	
[Synechocystis sp.]	2e-31
43 protein [Synechococcus sp.] 0.0	
OTOSYSTEM II D2 PROTEIN 76 pir JT0321 photosystem II protein D2 (M20815) D2 thylakoid protein sbDI thylakoid protein [Synechococcus]	
POTHEICAL 21.2 KD PROTEIN potheical protein 4 - Porphyra purpurea al chloroplast ORF 4, [Porphyra]	
TOLACETATE SYNTHASE SMALL CID SYNTHASE SMALL SUBUNIT acid synthase small chain - Porphyra cetohydroxyacid synthase small subunit	4.1.3.18 - map00290 map00650 map00660 map00770
trans isomerase B [Synechococcus]	



# Computational Biology at ORNL

Home About Us Analysis Tools Information Resources Projects & Research  
 Channel • Generation • Grail • GrailEXP • Pipeline • Parser • PROSPECT

Produced

for the [Joint](#)

[Genome](#)

[Institute](#)

Microbial

Sequencing

program

*Prochlorococcus*

## Modeled Genes Organization

- Metabolic Pathway
- Functional Category
- Gene
- Draft pmar Data - Netscape

Tab

- 
- 
- 
- 
- 
- 
- 

Search

OAK  
U. S. I

Prochlorococcus marinus files - Netscape  
 File Edit View Go Communicator Help

Version 10dec99 - Contig26 Gene 875

Gene Finders

Strand = r  
 Stop Location = 279721  
 Stop Codon = TAG

Gene Modeler	Start Location	Start Codon
Generation	279194	ATG
Glimmer	279194	ATG
Critica	279194	ATG

Blast against NR

Best Hit

gi|1669716|emb|CAA93118| (Z68890) class III phycoerythrin beta-subunit [Prochlorococcus marinus] >gi|5734498|emb|CAB52705.1| (AJ001230) phycoerythrin type III beta subunit [Prochlorococcus marinus]

Refresh BLAST results

PFam Model Comparison

Model	Description	Evalue	Score
Phycobosome	Phycobiosome proteins	9.4e-29	103.1

Refresh PFAM results

Best NR Hit with an EC Number

No EC Pathway hits

Synechocystis PCC6803 Comparison

Synecho ID	Description	Evalue	Percent Identity	Category
/	phycocyanin b subunit	4e-18	31.617647	

MRNA

```
ATGACAGTTCAAAAGAGTAATCAAATTATCAATGATAGAGATTAGAAAATAAAGTAATAAAAATATTGAAGACAT
AAAAGAACATTAAATAACTGCAAACTCAAGATTAGATGCAATAGATTCAAAACAATAATTCTCACGCAATTGCCGTG
ATGCTGTGACTGCAATGATTGTGAAAATCAAGATTGATTCAGTTAATCACAAAATATCTTGTAGATCACCACAAATAAGATGTCT
GTTTGTCTAAGAGATGGAGAAATAATTAAAGGATTGTTCTTACCTTTGATTTCTGATGACGAATCAGTTTATCTAA
AAACTGTTAAAGGATCTAAATAACTTATTTGCCCTTGGGTACCTCTGAAAAATGCTATCCGAGTTTTGAATTGA
TGAGAGATGCAACGATTCTGATTTAAAGTCACTGTAATCTGAAAGGAGAAAAAGAATTCTTCTGATTTAATT
TCTAATACAGAGTTCAATTGAGAGAATAATTAACTTTAAGATAG
```

PROTEIN

```
MTVSKSNQILSNRDLENISNKNIEDIKEFINTANSRLDAIDSITNNSHIAAADAVTAMICENQDSVNTKISLDTTNKMS
VCLRDGEIILRIVSYLLISDDESVLSKNCLKDLKNTYLALGVPLKNAIRVFELMRDATISDLKSTVNNSMKGEKEFLSDL
SNTEOFERIINLLR*
```

Document Done

Rhodopseudomonas palustris 1 - ORNL METABOLIC PATHWAYS

01100 5/100

Enzyme Maps

[Ec 4.1.1.18 - Synechocystis sp.]	
[Synechocystis sp.] 2e-31	
43 protein [Synechococcus sp.] 0.0	
[OSYSTEM II D2 PROTEIN 76 pir JT0321 photosystem II protein D2 (M20815) D2 thylakoid protein sbDI thylakoid protein [Synechococcus]	
[THEORETICAL 21.2 KD PROTEIN putative protein 4 - Porphyra purpurea al chloroplast ORF 4. [Porphyra]	
TOLACETATE SYNTHASE SMALL CID SYNTHASE SMALL SUBUNIT acid synthase small chain - Porphyra cetohydroxyacid synthase small subunit	4.1.3.18 - map00290 map00650 map00660 map00770
trans isomerase B [Synechococcus]	

## MED4 genes g

- P700 apoprotein subunit
- P700 apoprotein subunit
- photosystem I subunit V
- photosystem I subunit I
- photosystem I subunit
- chlorophyll a/b binding protein

D1 prot CP47 prot CP43 prot

D2 prot a subunit b subunit

base subunit

PsbH protein

PsbI protein

PsbJ protein

PsbK protein

PsbL protein

PsbM protein

PsbN protein

manganese

PsbT protein

13 kD protein

11 kD protein

erythrin

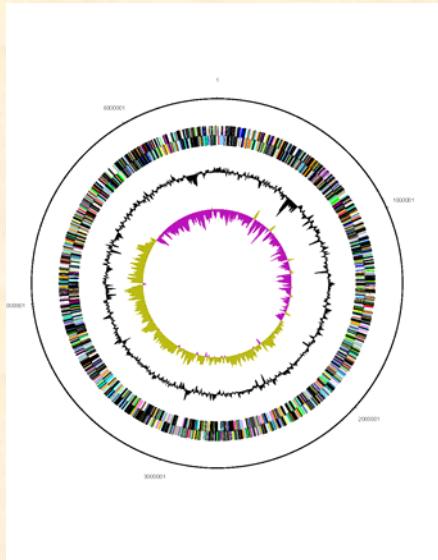
(cytM)

B)

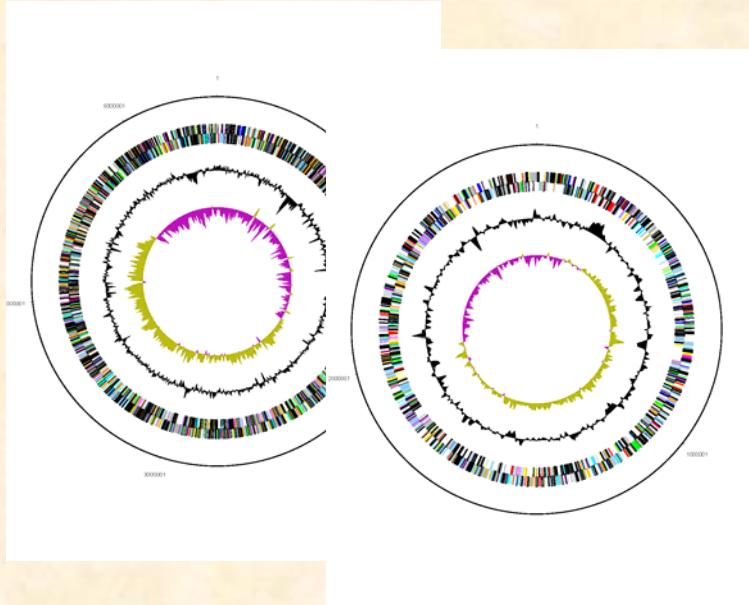
tE)

145

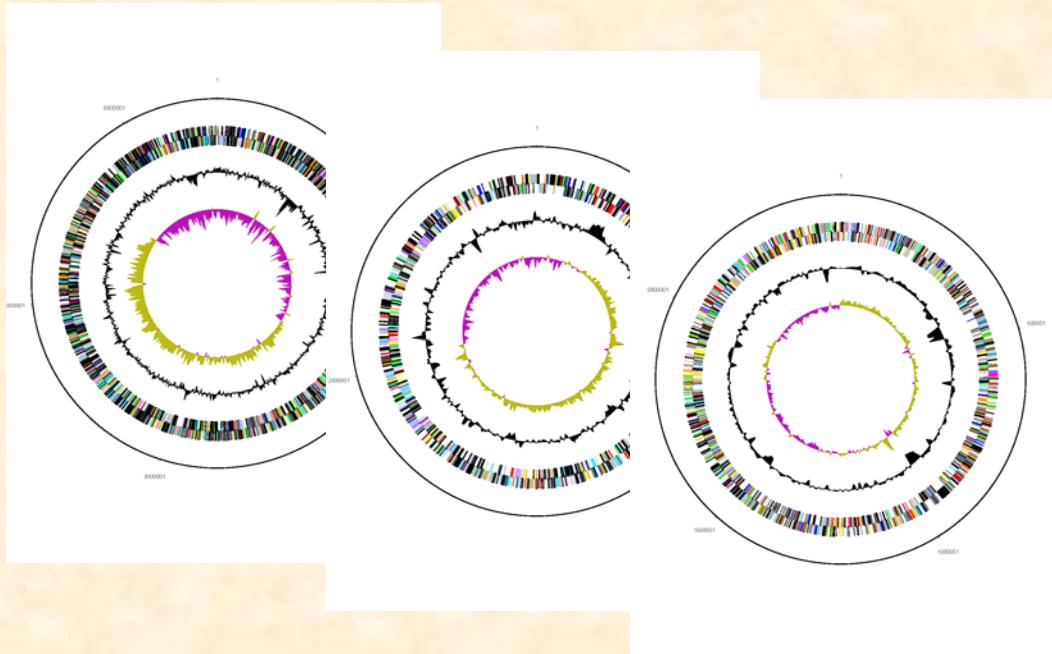
# Perspective has shifted from genes to genomes



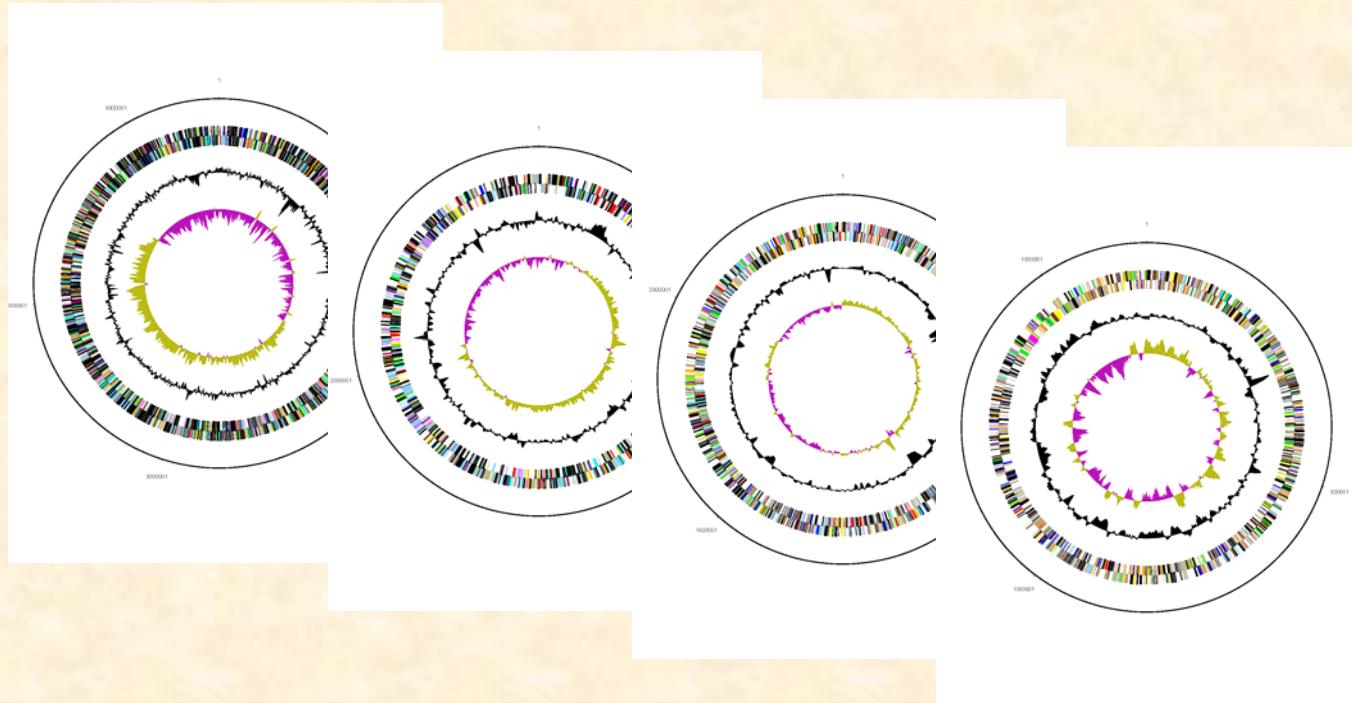
# Perspective has shifted from genes to genomes



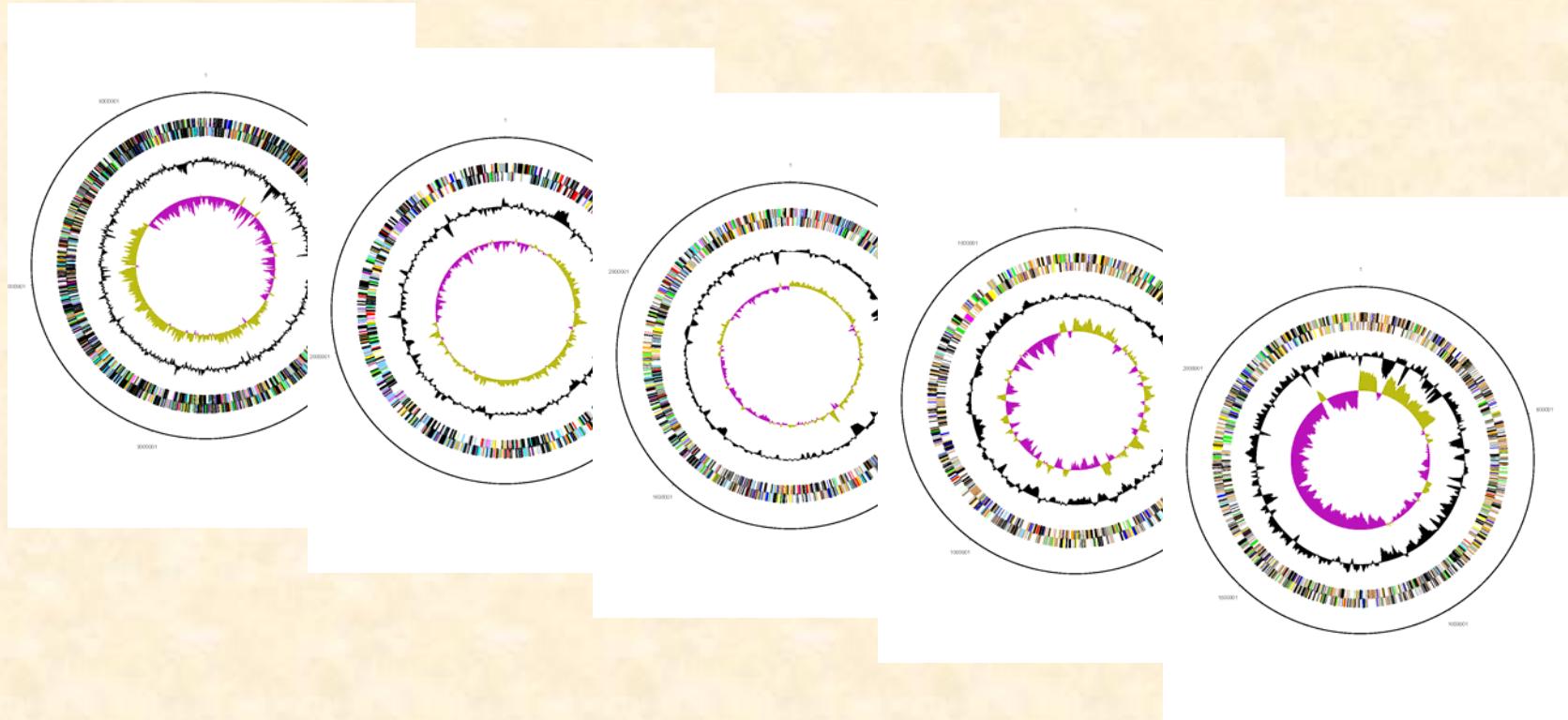
# Perspective has shifted from genes to genomes



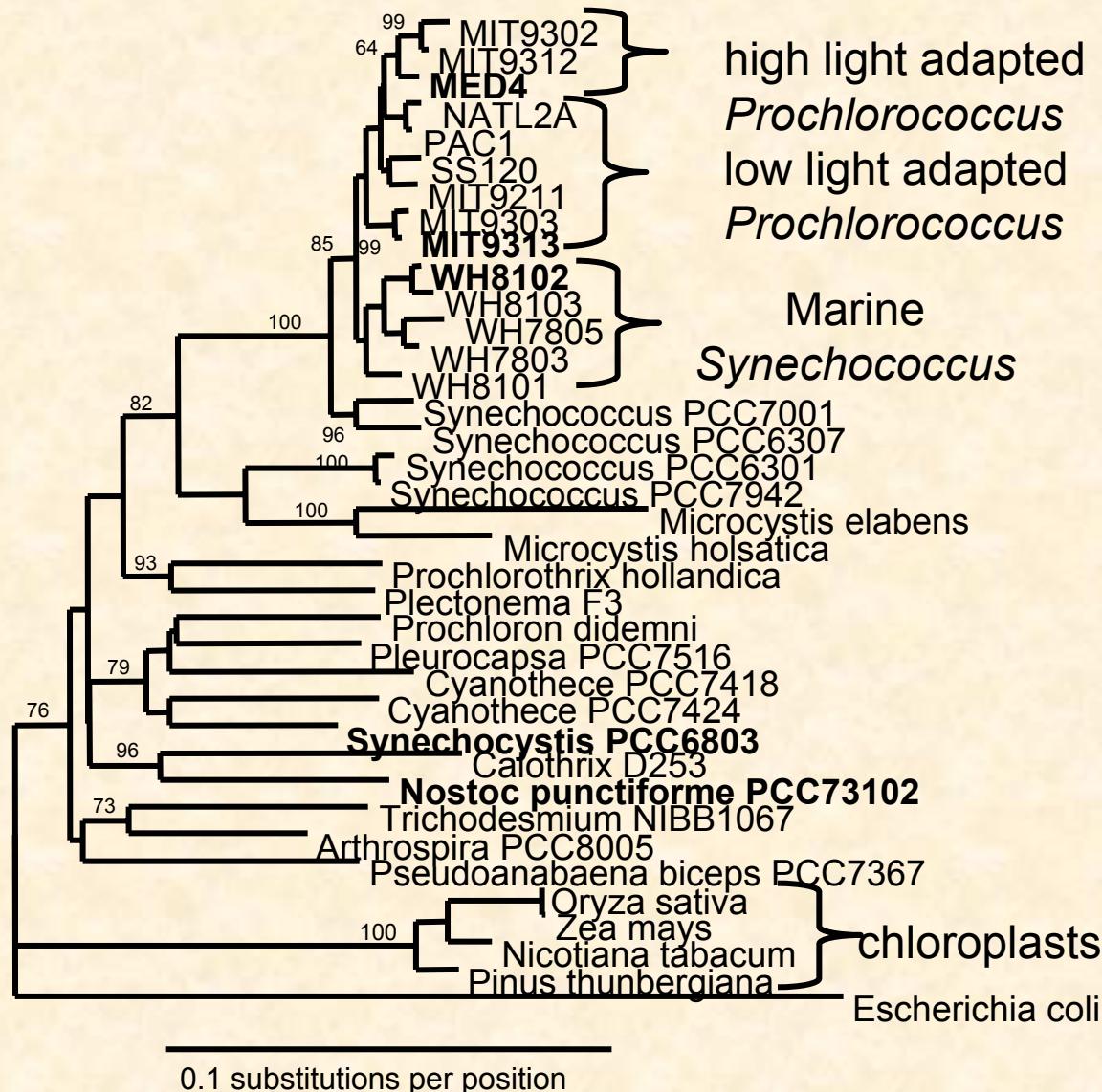
# Perspective has shifted from genes to genomes



# Perspective has shifted from genes to genomes

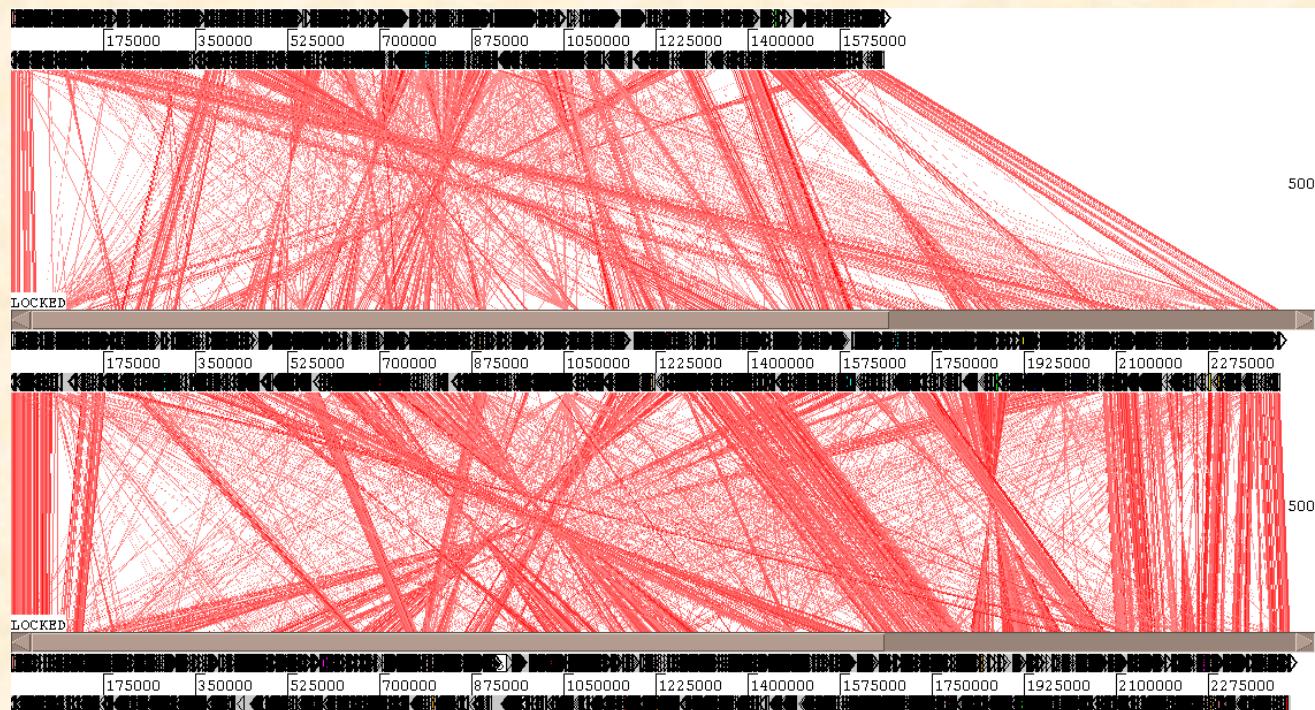


# Comparative genomics



# Comparative Genomics: Global alignment

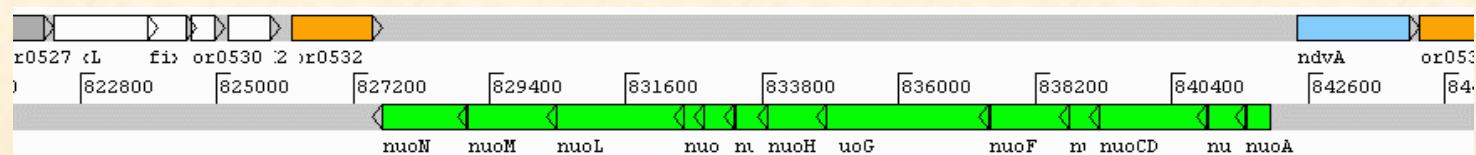
*P. marinus* MED4



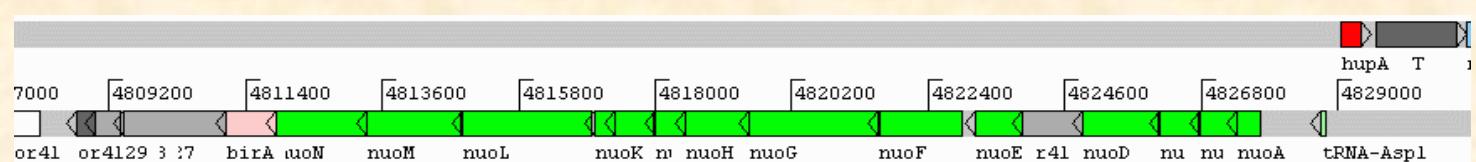
# NADH-ubiquinone oxidoreductase operon organization

## Unified operon

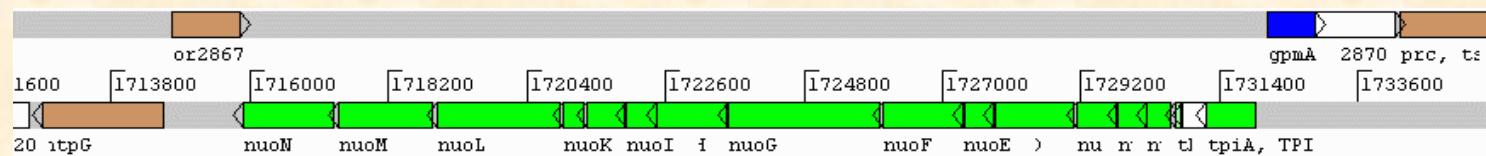
*R. palustris*



*R. palustris*

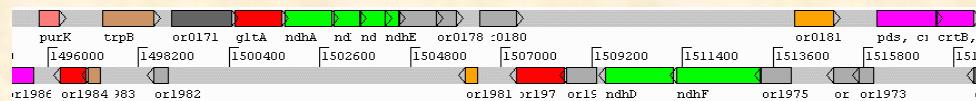


*N. europaea*

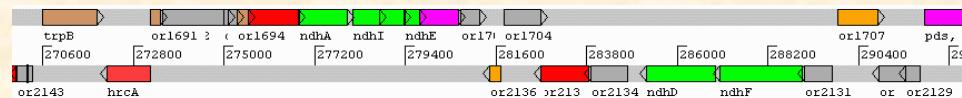


# NADH-ubiquinone oxidoreductase operon organization

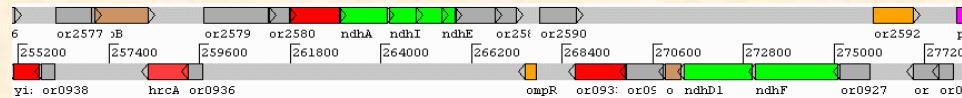
## Dispersed operons



*P. marinus* MED4

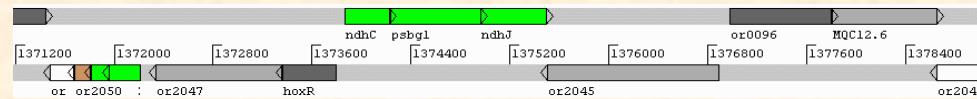


*P. marinus* MIT9313

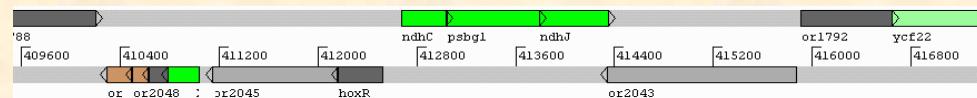


*Synechococcus* WH8102

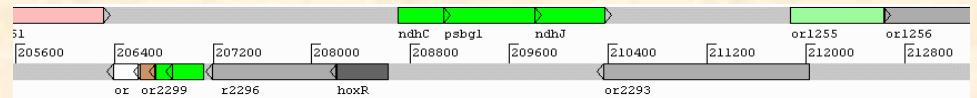
*P. marinus* MED4



*P. marinus* MIT9313

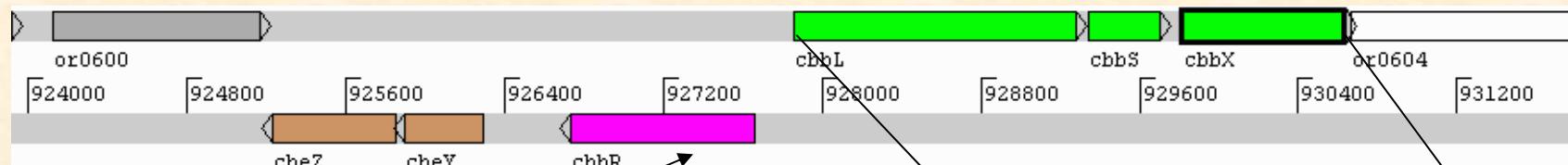


*Synechococcus* WH8102

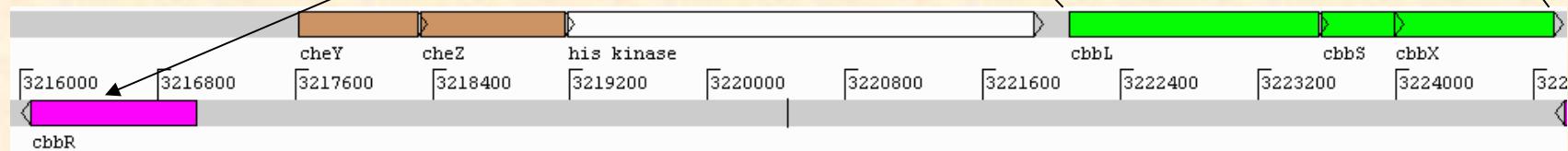


# Differing organization of RuBisCO operons

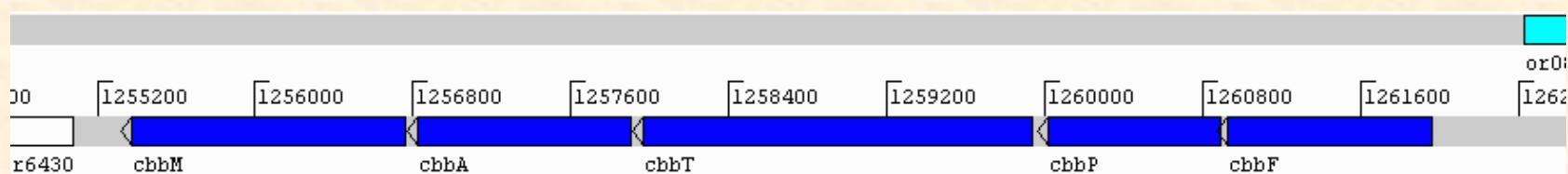
## *N. europaea*



## *R. palustris*

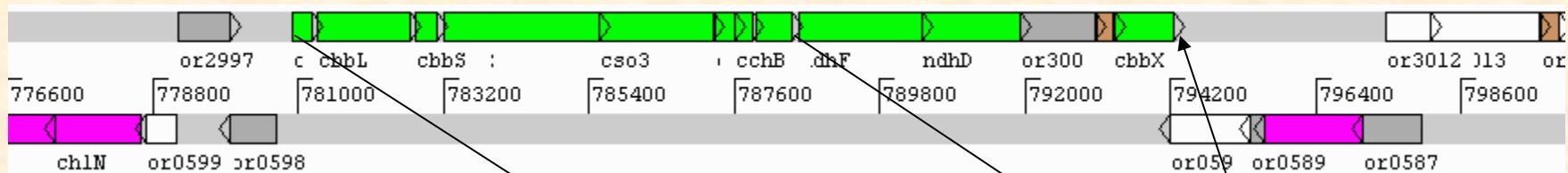


## *R. palustris*

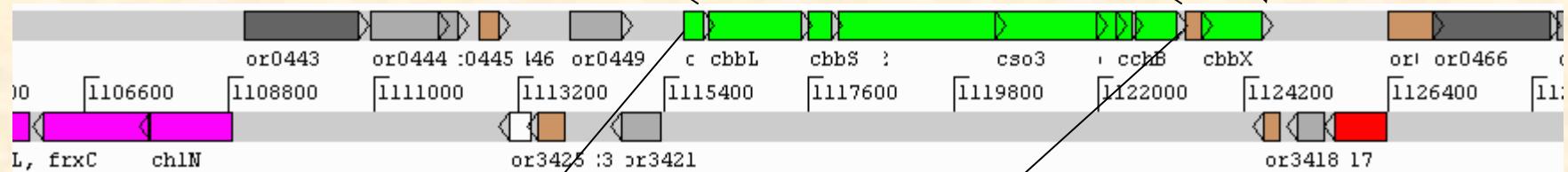


# Differing organization of RuBisCO operons

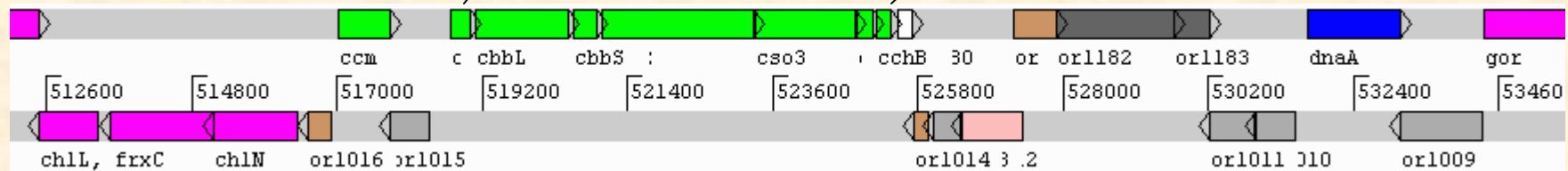
## *Synechococcus WH8102*



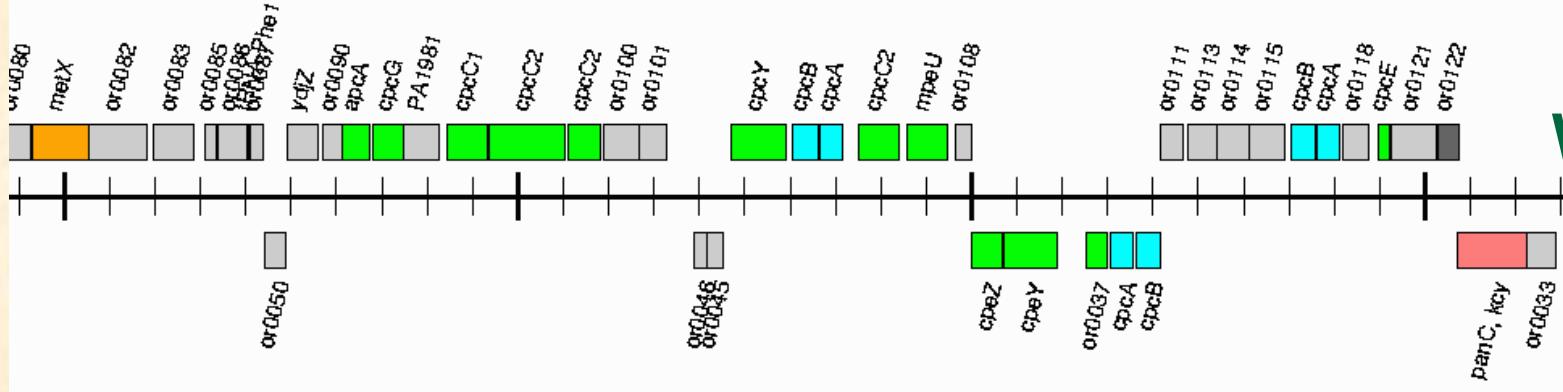
## *P. marinus MIT9313*



## *P. marinus MED4*



# Phycoerythrin and phycocyanin gene clusters



# Futures

- High-throughput microbial genomics
- White rot fungus
- Finished human chr5, 16, and 19
- Poplar
- Comparative genomics
  - genome evolution
  - operons, regulons, and control networks
  - proteomics
  - community perspective - “genomes of ecosystems”

# Acknowledgements

## ORNL Genome Analysis

Erich Baker  
Gwo-Liang Chen  
Michael Galloway  
Loren Hauser  
Douglas Hyatt  
Miriam Land  
Phil Locascio  
Victor Olman  
Denise Schmoyer  
Manesh Shah  
Jay Snoddy  
Inna Volker  
Ed Uberbacher

## Collaborators

Dan Arp  
Ron Atlas  
Penny Chisholm  
Dan Cullen  
Carrie Harwood  
Bernard Henrissat  
Mike Himmel  
Allen Hooper  
Martin Klotz  
Mark McBride  
Ron Mackenzie  
Jack Meeks  
Brian Palenik  
Gabrielle Rocap  
Bob Tabita

## JGI

Dan Rokhsar  
Stephanie Stilwagen  
Art Kobayashi  
Paul Predki

## LLNL

Patrick Chain  
Jane Lamerdin

## USDOE

Marv Frazier  
Dan Drell  
Anna Palmisano

# **Genome annotation: Bioinformatics for high- throughput genomics and beyond**

**Frank Larimer**  
**Oak Ridge National Laboratory**

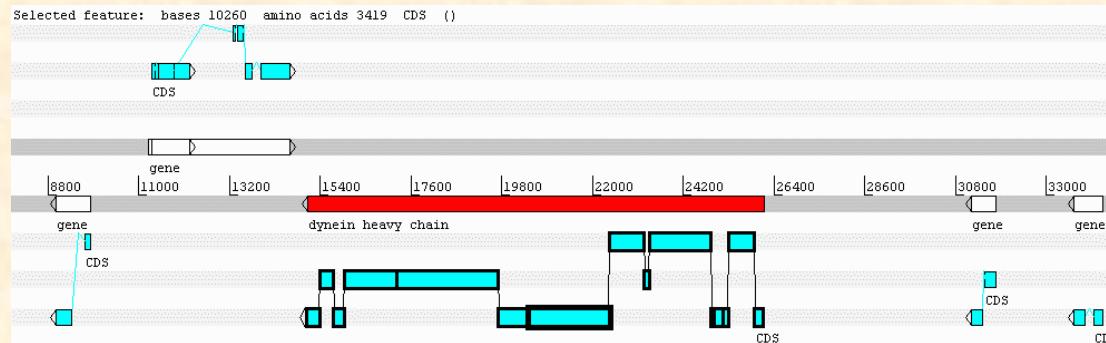
**9th DOE Contractor and Grantee Workshop**  
**Oakland, CA, 28 January 2002**

# **GrailEXP development for annotation of the white rot genome**

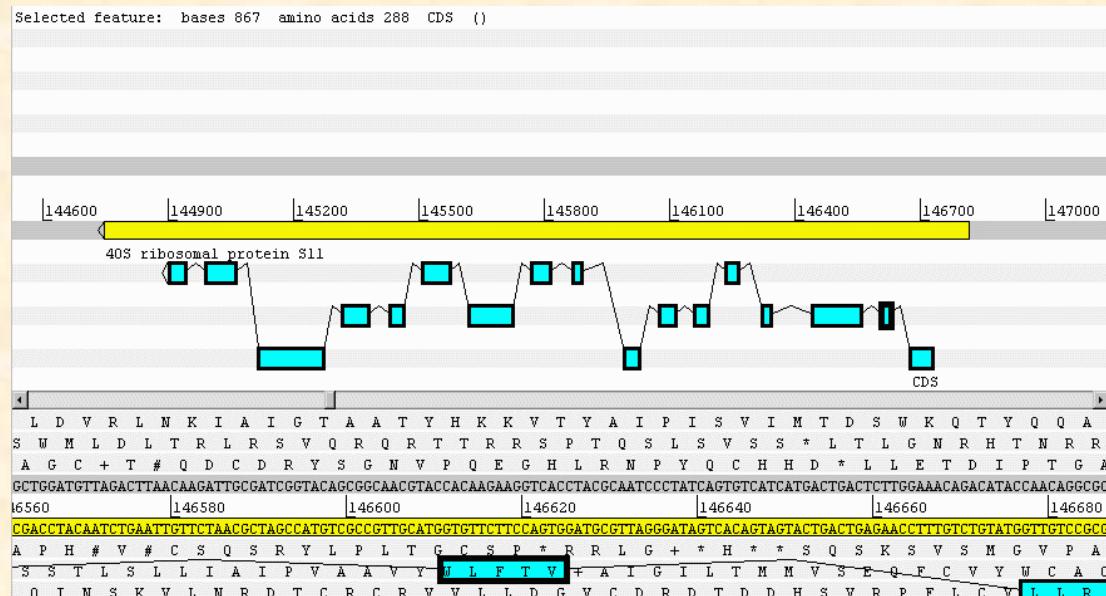
- A set of ~1300 white rot ESTs and mRNAs was aligned with the genomic sequence using GrailEXP
- Using the information from these ESTs, splice site scoring systems were constructed.
- GrailEXP was run on the genomic sequence using the custom white rot splice site scoring system and revised rules for intron sizes (rewarding introns close to the preferred size of ~50 bases).
- Iterative runs, with retraining, incorporating (prototype) protein alignment and heterologous ESTs

# Gene modeling: white rot

## Dynein: large exons



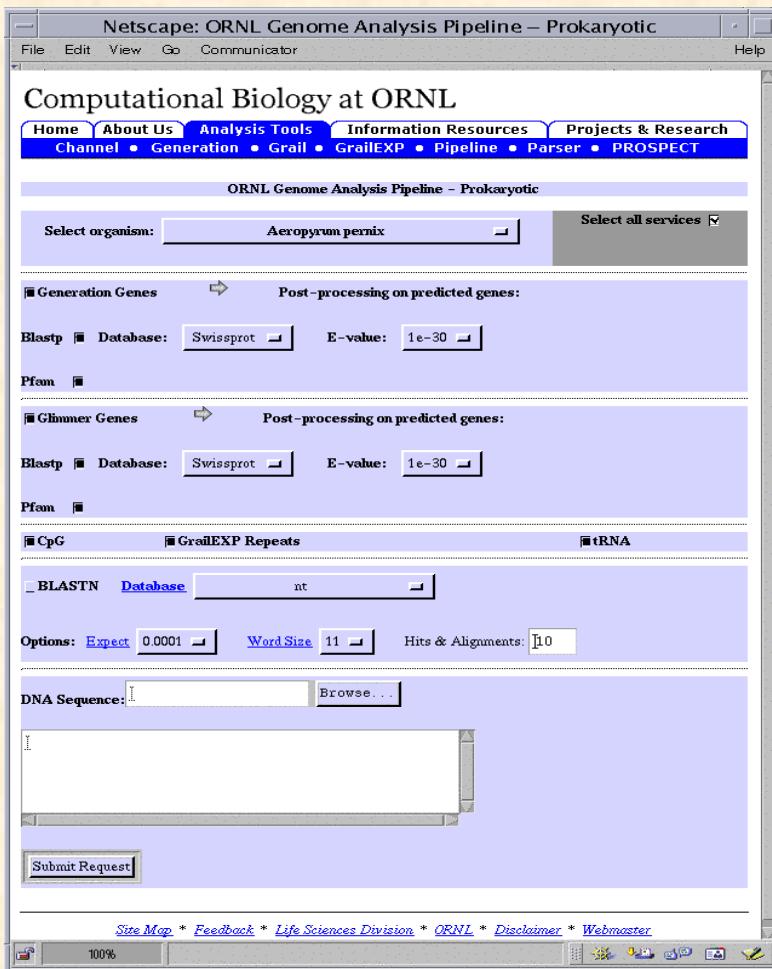
## S11: small exons



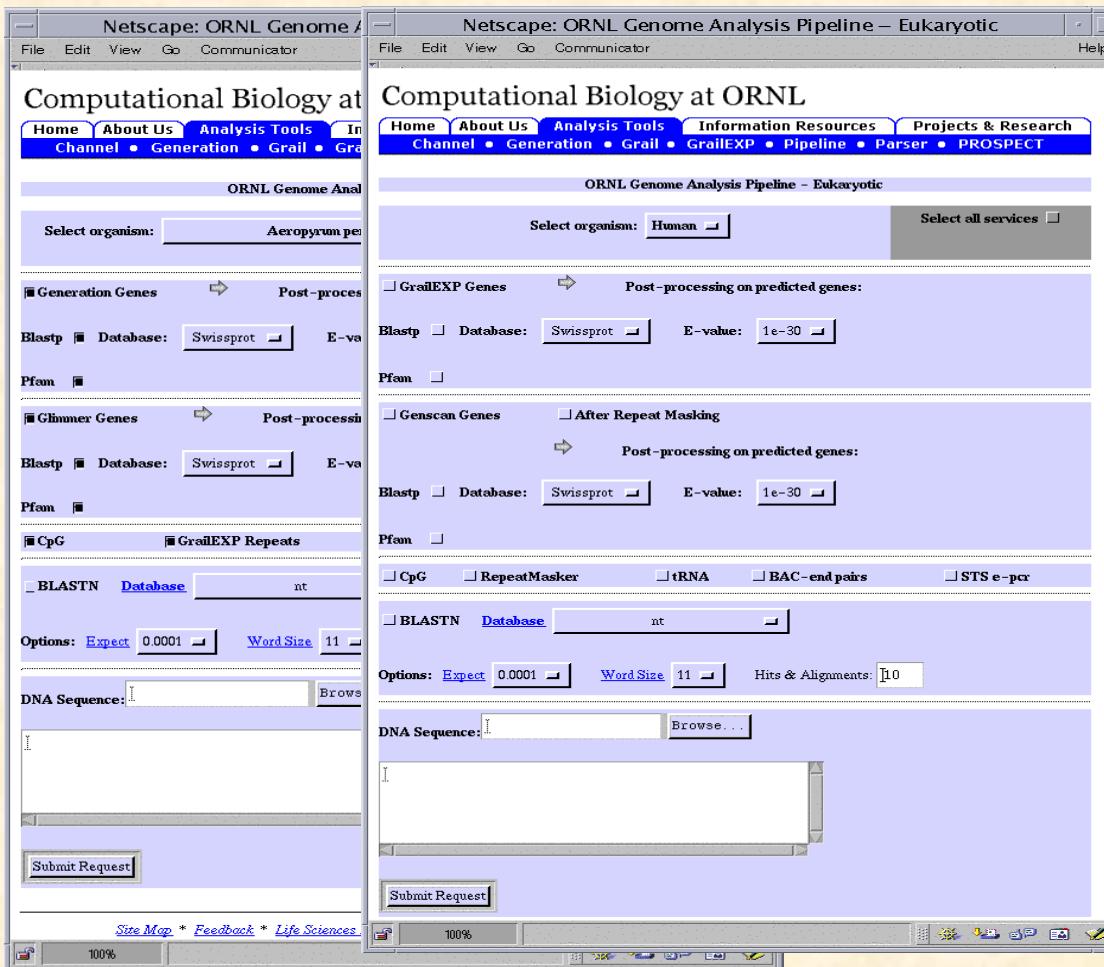
# **GrailEXP development**

- Training of GrailEXP on new genomes
- Streamlining and automation of the training process
- Recognition of more forms of alternative splicing
- Addition of protein homology to the system
- Addition of a TBLASTX-based approach for utilizing other genomic sequences, i.e. *Fugu* alignments to locate human genes
- Improvement of user interfaces (JAVA/HTML)
- Support for the DAS (Distributed Annotation System) protocol

# Analysis tools in action



# Analysis tools in action



# Analysis tools in action

The figure displays three side-by-side screenshots of the ORNL Genome Analysis Pipeline interface, showing the progression of a sequence analysis request through various tools.

**Screenshot 1 (Left):** Shows the initial input stage. A DNA sequence is pasted into a text area labeled "DNA Sequence:". Below it is a "Submit Request" button. At the bottom, there are links to "Site Map", "Feedback", and "Life Sciences Division".

**Screenshot 2 (Middle):** Shows the pipeline status after submission. It displays the following steps:

- GrailEXP Genes:** Status: Succeeded
- Proteins BLASTP:** Status: Succeeded
- Proteins Pfam:** Status: Succeeded
- Genscan Genes:** Status: Processing

At the bottom, there are links to "Site Map", "Feedback", "Life Sciences Division", "ORNL", "Disclaimer", and "Webmaster".

**Screenshot 3 (Right):** Shows the final output stage. It displays the results of the analysis:

- GrailEXP Genes:** Status: Succeeded
- Proteins BLASTP:** Status: Succeeded
- Proteins Pfam:** Status: Succeeded

At the bottom, there are links to "Site Map", "Feedback", "Life Sciences Division", "ORNL", "Disclaimer", and "Webmaster".

# Analysis tools in action

**Computational Biology at ORNL**

**Computational Biology at ORNL**

**Computational Biology at ORNL**

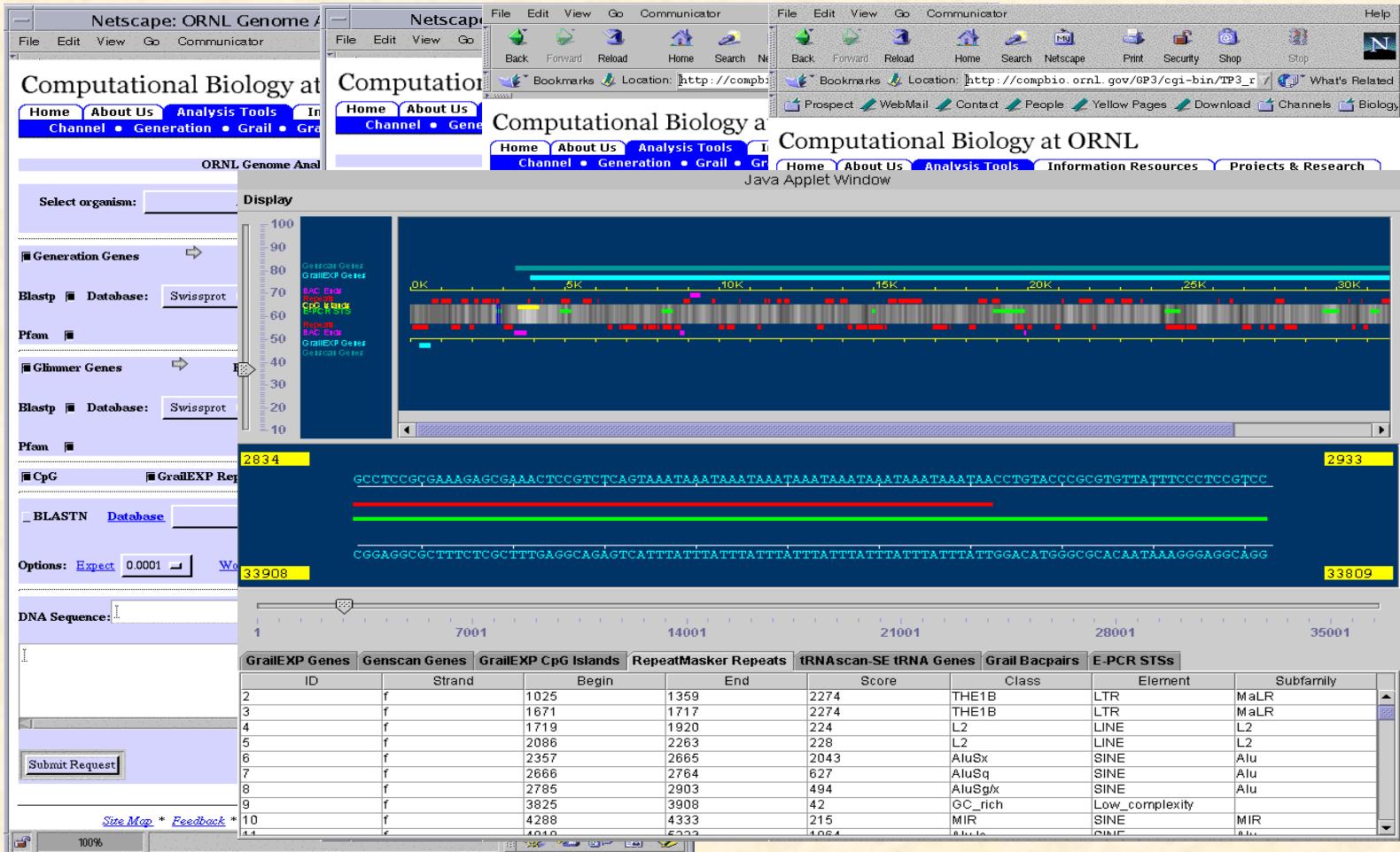
**Computational Biology at ORNL**

Organism Type: euk  
Organism: human  
Sequence header: >humadag  
Request size: 41150 bytes  
Sequence length: 36741 bp

Retrieval status: 3 found

Tool	Count	Status	Action
GrailEXP Genes	3	Succeeded	Retrieve
Proteins BLASTP	3	Succeeded	GeneID
Proteins Pfam	3	Succeeded	GeneID
Genscan Genes	1	Succeeded	Retrieve
Proteins BLASTP	1	Succeeded	GeneID
Proteins Pfam	1	Succeeded	GeneID
GrailEXP CpG Islands	1	Succeeded	Retrieve
RepeatMasker Repeats	88	Succeeded	Retrieve
tRNAscan-SE tRNA Genes	0	Succeeded	Retrieve
Grail BAC Pairs	4	Succeeded	Retrieve
E-PCR STSs	13	Succeeded	Retrieve

# Analysis tools in action



## Organization of Phycoerythrin (*cpe*) and Phycocyanin (*cpc*) Clusters

- *P. marinus* MED4
  - *cpeB* only
- *P. marinus* MIT9313
  - *cpeA cpeB*, seven others in conserved cluster
- *Synechococcus*
  - long *cpeA cpeB* cluster, multiple AB sets, inverted repeat, plus 13 others
- *Nostoc*
  - 3 sets of *cpcB cpcA* clusters
- *Synechocystis*
  - 1 *cpcA cpcB*, 1 *apcA apcB*, solo *apcB*

# *kaiABC* clock operon

